

# SHADE: INFORMATION-BASED REGULARIZATION FOR DEEP LEARNING

Michael Blot<sup>1</sup>, Thomas Robert<sup>1</sup>, Nicolas Thome<sup>2</sup>, Matthieu Cord<sup>1</sup>

<sup>1</sup> Sorbonne Universités, UPMC Univ Paris 06, CNRS, LIP6 UMR 7606, 4 place Jussieu, 75005 Paris

<sup>2</sup> CEDRIC, Conservatoire National des Arts et Métiers, 292 Rue St Martin, 75003 Paris, France

{michael.blot, thomas.robert, matthieu.cord}@lip6.fr - nicolas.thome@cnam.fr

## ABSTRACT

Regularization is a big issue for training deep neural networks. In this paper, we propose a new information-theory-based regularization scheme named SHADE for SHAnnon DEcay. The originality of the approach is to define a prior based on conditional entropy, which explicitly decouples the learning of invariant representations in the regularizer and the learning of correlations between inputs and labels in the data fitting term. Our second contribution is to derive a stochastic version of the regularizer compatible with deep learning, resulting in a tractable training scheme. We empirically validate the efficiency of our approach to improve classification performances compared to standard regularization schemes on several standard architectures.

**Index Terms**— Deep learning, regularization, invariance, information theory, image understanding

## 1. INTRODUCTION

Deep neural networks (DNNs) have shown impressive state-of-the-art results in the last years on numerous tasks and especially for image classification [1, 2]. One key element is the use of very deep models with a huge number of parameters. Accordingly, DNNs need to be trained on a lot of data (*e.g.* ImageNet) and with a regularization scheme to control overfitting. Although regularization methods such as weight decay [3], dropout [4] or batch normalization [5] are common practice, the question of DNN regularization remains open as demonstrated by [6].

Formally, let us note  $X \in \mathcal{X}$  the input variable,  $C \in \mathcal{C}$  the output (class) variable,  $w$  model parameters and  $Y = h(w, X)$  the (deep) representation of the input that leads to the class prediction. Usually, training schemes for deep models for classification tasks use an objective function which linearly combines a classification loss  $\ell_{\text{cls}}(w, X, Y, C)$  – generally cross-entropy – and a regularization term  $\Omega(w, X, Y, C)$ , with  $\beta \in \mathbb{R}^+$ :

$$\mathcal{L}(w) = \mathbb{E}_{(X,C)}(\ell_{\text{cls}}(w, X, Y, C) + \beta \cdot \Omega(w, X, Y, C)) \quad (1)$$

This paper studies the issue of regularization, and we propose a new regularization term  $\Omega(w, X, Y, C)$  in Eq (1).

**Quantifying invariance.** Designing DNN models that are robust to variations on the training data and that preserve class information is the main motivation of this work. With this motivation, Scattering decompositions [7] are appealing transforms, which have been incorporated into adapted network architectures like [8]. However, for tasks like image recognition, it is very difficult to design an explicit modelling of all transformations a model should be invariant to.

**Information-theory-based regularization.** Many works like [9] use information measures as regularization criterion. The Information Bottleneck framework (IB) proposed in [10] suggests to use mutual information  $I(X, Y)$  (see [11] for definition) as regularization criterion. [12] extends it to a variational context VIB. However, regularization based on  $I(X, Y)$  may conflict with the task loss  $\ell_{\text{cls}}(w, X, C)$  in Eq (1). In addition, IB [10] is computationally expensive and is only applied at a single final layer of the network.

In this paper, we propose a new regularization method, denoted as SHADE for SHAnnon DEcay. Our first contribution is to design a new regularization loss, that aims at minimizing a particular criterion: the entropy of the representation variable conditionally to the class variable, *i.e.*  $\Omega(w, X, C) = H(Y | C)$ . This criterion strongly supports intra-class invariance of the representation, without conflicting with  $\ell_{\text{cls}}(w, X, C)$  in Eq (1). Our second contribution consists in deriving a tractable surrogate function of  $H(Y | C)$ . This enables the incorporation of the regularizer at every layer of the network, leading to a scalable optimization scheme based on stochastic gradient descent (SGD).

We provide an extensive experimental validation of our SHADE regularizer for important standard DNNs, namely AlexNet, ResNet and Inception, applied to CIFAR-10 and ImageNet datasets.

## 2. SHADE: A NEW REGULARIZATION METHOD

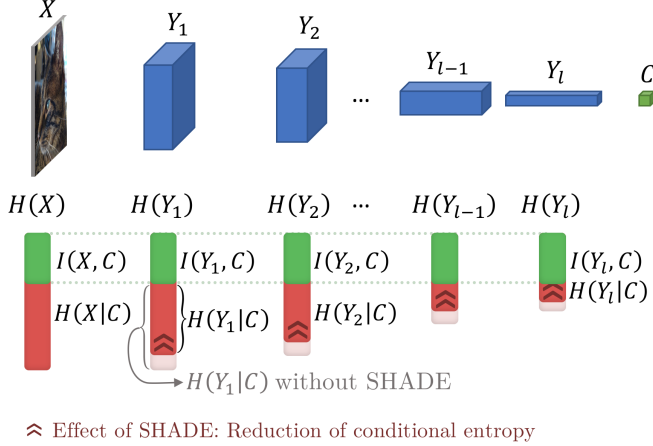
In this section, we further describe SHADE, a new regularization term based on the conditional entropy  $H(Y | C)$  designed to drive the optimization towards a more invariant representation.

### 2.1. Conditional Entropy-based Regularization for Deep Neural Networks

In this article, the considered task is the classification of images, so we focus on intra-class invariance, explaining the use of  $H(Y | C)$  as a criterion. An overview of the approach is given in Fig. 1.

Our approach differs from the Information Bottleneck framework (IB) [10], which suggests to use  $I(X, Y)$  as a regularizer. In the case where  $\mathcal{X}$  is discrete, our criterion is related to IB's through the following development  $I(X, Y) = I(C, Y) + H(Y | C)$  that holds for deterministic models ( $H(Y | X) = 0$ ). In a context of optimization with SGD, minimizing  $H(Y | C)$  appears to be more efficient to preserve the term  $I(C, Y)$ , which represents the mutual information of the representation variable with the class variable and must stay high to predict accurately  $C$  from  $Y$ . It is illustrated in Fig. 1 where we see that  $I(Y_i, C)$  (in green) remains constant while  $H(Y_i | C)$  (in red) decreases.

We claim that  $H(Y | C)$  quantifies accurately how invariant a representation is, while being agnostic to the transformations it is invariant to. When developing the entropy of the representation,



**Fig. 1:** DNN architecture with corresponding layers’ entropies, showing the layer-wise action of SHADE. Given that  $H(Y_i) = I(Y_i, C) + H(Y_i | C)$ , SHADE minimizes  $H(Y_i | C)$  without affecting  $I(Y_i, C)$ .

we get  $H(Y) = I(X, Y) + H(Y | X)$ . Given that the DNN is deterministic, we have  $H(Y | X) = 0$  so  $H(Y) = I(X, Y) = H(X) - H(X | Y)$ .  $H(X)$  being fixed,  $H(Y)$  is inversely related to  $H(X | Y)$ . This last term can lower bound any input reconstruction error, demonstrated with inequalities such as Fano’s inequality [11]. Thus, it can perfectly quantify how difficult it is to recover the input from its representation. The benefit of compressing representations is also consistent with Occam’s Razor interpretation of the “minimum description length” principle of [13].

**Layer-wise regularization.** A DNN is composed of a number  $L$  of layers that transform sequentially the input. Each one can be seen as an intermediate representation variable, noted  $Y_\ell$  for layer  $\ell$ , that is determined by the output of the previous layer and a set of parameters  $w_\ell$ . Each layer filters out part of the information from the initial input. Thus, from the data processing inequality in [11], the following inequalities can be derived for any layer  $\ell$ :

$$H(Y_\ell | C) \leq H(Y_{\ell-1} | C) \leq \dots \leq H(Y_1 | C) \leq H(X | C).$$

The conditional entropy of a layer impacts the conditional entropy of the subsequent layers. Consequently, similarly to the recommendation of [14], as illustrated on Fig. 1, we apply a regularization on all layers, minimizing the layer-wise entropy  $H(Y_\ell | C)$ , and producing a global criterion:

$$\Omega_{\text{layers}} = \sum_{\ell=1}^L H(Y_\ell | C). \quad (2)$$

**Unit-wise regularization.** Examining one layer  $\ell$ , its representation variable is a random vector of coordinates  $Y_{\ell,i}$  and of dimension  $D_\ell$ :  $Y_\ell = (Y_{\ell,1}, \dots, Y_{\ell,D_\ell})$ . The upper bound<sup>1</sup>  $H(Y_\ell | C) \leq \sum_{i=1}^{D_\ell} H(Y_{\ell,i} | C)$  allows us to consider the different units of a layer independently and then to define a unit-wise criterion that SHADE seeks to minimize. For each unit  $i$  of every layer  $\ell$  we design a

<sup>1</sup>This upper bound is well justified in deep learning as the neurons of a layer tend to be more and more independent of each other as we go deeper within the network.

loss  $\omega_{\text{unit}}(Y_{\ell,i} | C) = H(Y_{\ell,i} | C)$  that will be part of the global regularization loss:

$$\Omega_{\text{layers}} \leq \Omega_{\text{units}} = \sum_{\ell=1}^L \sum_{i=1}^{D_\ell} \underbrace{H(Y_{\ell,i} | C)}_{\omega_{\text{unit}}(Y_{\ell,i} | C)}. \quad (3)$$

In the sections that follow, we use the notation  $Y$  instead of  $Y_{\ell,i}$  for simplicity since the coordinates are all considered independently to define  $\omega_{\text{unit}}(Y_{\ell,i} | C)$ .

## 2.2. Estimating Entropy

In this section, we describe how to define a loss based on the measure  $H(Y | C)$  with  $Y$  being one coordinate variable of one layer. Defining this loss is not obvious as the gradient of  $H(Y | C)$  with respect to the layer’s parameters may be computationally intractable.  $Y$  has an unknown distribution and without modeling it properly it is not possible to compute  $H(Y | C)$  precisely for the following reasons.

Since  $H(Y | C) = \sum_{c \in \mathcal{C}} p(c) H(Y | c)$  most of the estimators require to compute  $|\mathcal{C}|$  different entropies  $H(Y | c)$ . This means that, given a batch, the number of samples used to estimate one of these entropies is divided by  $|\mathcal{C}|$  on average which becomes particularly problematic when dealing with a large number of classes such as the 1,000 classes of ImageNet. Furthermore, entropy estimators are extremely inaccurate considering the number of samples in a batch. For example, MLE estimators of entropy in [15] converge in  $\mathcal{O}((\log K)^2/K)$  for  $K$  samples. Finally, most estimators such as MLE require discretizing the space in order to approximate the distribution *via* a histogram. This raises issues on the bins definition considering that the variable distribution is unknown and varies during the training in addition to the fact that having a histogram for each neuron is computationally and memory consuming.

To tackle these drawbacks we investigate the two following tricks: the introduction of a latent variable that enables to use more examples to estimate the entropy; and a bound on the entropy of the variable by an increasing function of its variance to avoid the issue of entropy estimation with a histogram, making the computation tractable and scalable.

**Latent code  $Z$ .** Intermediate features of a DNN most likely take similar values for inputs of different classes – this is especially true for low-level features. The semantic information provided by a single feature  $Y$  therefore describes a particular pattern it detects rather than the detection of a class. Only the association of features allows determining the class. To better understand how the class information is encoded in an individual neuron variable (before ReLU), let us take a look at the behavior of the activation function used. The ReLU activation makes the neuron act as a detector, returning a signal when a certain pattern is present on the input. If the pattern is absent the signal is zero, otherwise,  $Y$  quantifies the resemblance with it.

We therefore propose to associate a Bernoulli variable  $Z$  to each unit variable  $Y$ . This variable  $Z$  indicates if a particular pattern is present on the input ( $Z = 1$  when  $Y \gg 0$ ) or not ( $Z = 0$  when  $Y \leq 0$ ). It acts like a latent code in variational models [16] or in generative models [17].

In other words,  $Z$  is a semantically meaningful factor about the class  $C$  and from which the input  $X$  is generated. The feature value  $Y$  is then a quantification of the possibility for this attribute to be present ( $Z = 1$ ) or not ( $Z = 0$ ) in the input. For instance, for high level features,  $Z$  could represent the presence or not of a particular object, that allows to discriminate between classes (e.g. for a wheel,

---

**Algorithm 1** Moving average updates: for  $z \in \{0, 1\}$ ,  $p^z$  estimates  $p(Z = z)$  and  $\mu^z$  estimates  $\mathbb{E}(Y | Z = z)$

---

```

1: Initialize:  $\mu^0 = -1, \mu^1 = 1, p^0 = p^1 = 0.5, \lambda = 0.8$ 
2: for each mini-batch  $\{y^{(k)}, k \in 1..K\}$  do
3:   for  $z \in \{0, 1\}$  do
4:      $p^z \leftarrow \lambda p^z + (1 - \lambda) \frac{1}{K} \sum_{k=1}^K p(z | y^{(k)})$ 
5:      $\mu^z \leftarrow \lambda \mu^z + (1 - \lambda) \frac{1}{K} \sum_{k=1}^K \frac{p(z | y^{(k)})}{p^z} y^{(k)}$ 
6:   end for
7: end for

```

---

presents on cars and trucks,  $Y$  notifies on the resemblance with a certain pattern representing a wheel while  $Z$  indicates if the wheel is present or not on the image). Note that  $Z$  is not a deterministic variable of  $Y$ .

Therefore we assume the Markov chain  $C \rightarrow Z \rightarrow X \rightarrow Y$  (see definition in [11]). We indeed expect  $Y$  to evolve towards a sufficient statistic of  $Z$  for  $C$  during the training. Considering the sufficient statistic relation  $I(Y, C) = I(Y, Z)$  we get the equivalent equality  $H(Y | C) = H(Y | Z)$ , to finally obtain:

$$\begin{aligned} \omega_{\text{unit}}(Y | C) &= H(Y | C) = H(Y | Z) \\ &= \sum_{z \in \{0,1\}} p(z) H(Y | Z = z). \end{aligned}$$

This modeling of  $Z$  as a binomial variable (one for each unit) has the advantage of enabling good estimators of conditional entropy since we only divide the batch into two sets for the estimation ( $z = 0$  and  $z = 1$ ) regardless of the number of classes.

**Variance bound.** The previous trick allows computing fewer entropy estimates to obtain the global conditional entropy, thus increasing the sample size used for each entropy estimation. Unfortunately, it does not solve the bin definition issue. To address this, we propose to use the following bound on  $H(Y | Z)$ , that does not require the definition of bins:  $H(Y | Z) \leq \frac{1}{2} \ln(2\pi e \text{Var}(Y | Z))$ .

This bound holds for any continuous distributions  $Y$  and there is equality if the distribution is Gaussian, which is a proper law to model the activations, according to [18]. For many other distributions such as the exponential ones, the entropy is also directly equal to an increasing function of the variance. In addition, one main advantage is that variance estimators are much more robust than entropy estimators, converging in  $\mathcal{O}(1/K)$  for  $K$  samples instead of  $\mathcal{O}(\log(K)^2/K)$ .

Finally, the  $\ln$  function being one-to-one and increasing, we only keep the simpler term  $\text{Var}(Y | Z)$  to design our final loss:

$$\Omega_{\text{SHADE}} = \sum_{\ell=1}^L \sum_{i=1}^{D_\ell} \sum_{z \in \{0,1\}} p(Z_{\ell,i} = z | Y) \text{Var}(Y | Z_{\ell,i} = z).$$

In next section, we detail the definition of the differential loss using  $\text{Var}(Y | Z)$  as a criterion computed on a mini-batch.

### 2.3. Instantiating SHADE

For one unit of one layer, the previous criterion writes:

$$\begin{aligned} \text{Var}(Y | Z) &= \int_{\mathcal{Y}} p(y) \int_{\mathcal{Z}} p(z | y) (y - \mathbb{E}(Y | z))^2 dz dy \quad (4) \\ &\approx \frac{1}{K} \sum_{k=1}^K \left[ \int_{\mathcal{Z}} p(z | y^{(k)}) (y^{(k)} - \mathbb{E}(Y | z))^2 dz \right] \quad (5) \end{aligned}$$

The quantity  $\text{Var}(Y | Z)$  can be estimated with Monte-Carlo sampling on a mini-batch of input-target pairs  $\{(x^{(k)}, c^{(k)})\}_{1 \leq k \leq K}$  of intermediate representations  $\{y^{(k)}\}_{1 \leq k \leq K}$  as in Eq. (5).

$p(Z | y)$  can be interpreted as the probability of presence of attribute  $Z$  on the input, so it should clearly be modeled such that  $p(Z = 1 | y)$  increases with  $y$ . We suggest using:

$$p(Z = 1 | y) = \sigma(y) \quad p(Z = 0 | y) = 1 - \sigma(y)$$

with  $\sigma(y) = 1 - e^{-\text{ReLU}(Y)}$ .

For the expected values  $\mu^z = \mathbb{E}(Y | z)$  we use a classic moving average that is updated after each batch as described in Algorithm 1. Note that the expectations are not changed by the optimization since they have no influence on the entropy  $H(Y | Z)$ .

For this proposed instantiation, our SHADE regularization penalty takes the form:

$$\Omega_{\text{SHADE}} = \sum_{\ell=1}^L \sum_{i=1}^{D_\ell} \sum_{k=1}^K \sum_{z \in \{0,1\}} p(Z_{\ell,i} = z | y_{\ell,i}^{(k)}) (y_{\ell,i}^{(k)} - \mu_{\ell,i}^z)^2.$$

We have presented a regularizer that is applied layer-wise and that can be integrated into the usual optimization process of a DNN. The additional computation and memory usage induced by SHADE is almost negligible (computation and storage of two moving averages per neuron). Namely, SHADE adds half as many parameters as batch normalization does.

## 3. EXPERIMENTS

### 3.1. Image Classification with Various Architectures on CIFAR-10

**Table 1:** Classification accuracy (%) on CIFAR-10 test set.

	MLP	AlexNet	ResNet	Inception
No regul.	62.38	83.25	89.84	90.97
Weight decay	62.69	83.54	91.71	91.87
Dropout	65.37	85.95	89.94	91.11
SHADE	66.05	85.45	<b>92.15</b>	<b>93.28</b>
SHADE+D	<b>66.12</b>	<b>86.71</b>	92.03	92.51

We perform image classification on the CIFAR-10 dataset, which contains 50k training images and 10k test images of  $32 \times 32$  RGB pixels, fairly distributed within 10 classes [19]. Following the architectures used in [6], we use a small Inception model, a three-layer MLP, and an AlexNet-like model with 3 convolutional and 2 fully connected layers. We also use a ResNet architecture from [20]. Those architectures represent a large family of DNN and some have been well studied in [6] within the generalization scope. For training, we use randomly cropped images of size  $28 \times 28$  with random horizontal flips. For testing, we simply center-crop  $28 \times 28$  images. We use momentum SGD for optimization, same protocol as [6].

We compare SHADE with two regularization methods: *weight decay* and *dropout*. For all architectures, the regularization parameters have been cross-validated to find the best ones for each method and the obtained accuracies on the test set are reported in Table 1.

We obtain the same trends as [6], which provides a small improvement of 0.31% over weight decay on AlexNet. The improvement over weight decay is slightly more important with ResNet and Inception

(0.87% and 0.90%) probably thanks to the use of batch normalization. In our experiments dropout improves generalization performances only for AlexNet and MLP. It is known that the use of batch normalization lowers the benefit of dropout, which is in fact not used in [2].

We first notice that for all kinds of architectures the use of SHADE significantly improves the generalization performance. It demonstrates the ability of SHADE to regularize the training of deep architectures.

Finally, SHADE shows better performances than dropout on all architecture except AlexNet, for which they seem to be complementary, probably because of the very large number of parameters in the fully-connected layers, with best performances obtained with SHADE coupled with dropout (named SHADE+D). This association is also beneficial for MLP. On Inception and ResNet, even if dropout and SHADE independently improve generalization performances, their association is not as good as SHADE alone, probably because it enforces too much regularization.

### 3.2. Large Scale Classification on ImageNet

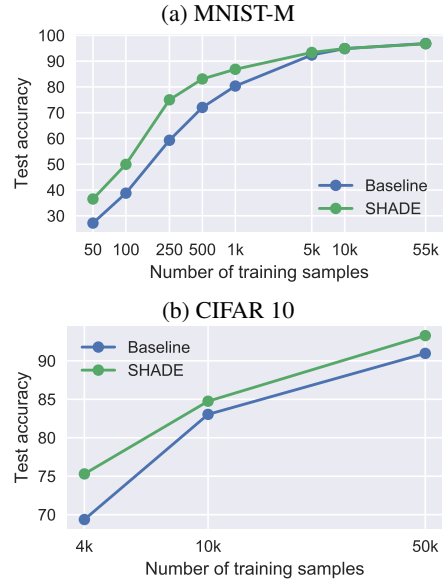
In order to test SHADE regularization on a very large scale dataset, we train on ImageNet [21] a WELDON network from [22] adapted from ResNet-101. This architecture changes the forward and pooling strategy by using the network in a fully-convolutional way and adding a max+min pooling, thus improving the performance of the baseline network. We used the **pre-trained weights of ResNet-101** (from the torchvision package of PyTorch) achieving performances on the test set of **77.56%** for top-1 accuracy and 93.89% for top-5 accuracy. Provided with the pre-trained weights, the **WELDON architecture** obtains **78.51%** for top-1 accuracy and 94.65% for top-5 accuracy. After fine tuning the network using **SHADE** for regularization we finally obtained **80.14%** for top-1 accuracy and 95.35% for top-5 accuracy for a concrete improvement. This demonstrates the ability to apply SHADE on very large scale image classification successfully.

### 3.3. Training with a Limited Number of Samples

When datasets are small, DNNs tend to overfit quickly and regularization becomes essential. Because it tends to filter out information and make the network more invariant, SHADE seems to be well fitted for this task. To investigate this, we propose to train DNNs with and without SHADE on MNIST-M [23] with different numbers of samples in the training set.

First, we tested this approach on the digits dataset MNIST-M. This dataset consists of the MNIST digits where the background and digit have been replaced by colored and textured information (see Fig. 3 for examples). The interest of this dataset is that it contains lots of unnecessary information that should be filtered out, and is therefore well adapted to measure the effect of SHADE. A simple convolutional network has been trained with different numbers of samples of MNIST-M and the optimal regularization weight for SHADE has been determined on the validation set. The results can be seen on Figure 2a. We can see that especially for small numbers of training samples (< 1000), SHADE provides an important gain of 10 to 15% points over the baseline. This shows that SHADE helped the model in finding invariant and discriminative patterns using less data samples.

Additionally, Figure 3 shows samples that are misclassified by the baseline model but correctly classified when using SHADE. These images contain a large amount of intra-class variance (color, texture, etc.) that is not useful for the classification task, explaining why



**Fig. 2:** Results when training with a limited number of samples with and without SHADE for MNIST-M (a), and CIFAR 10 (b).



**Fig. 3:** Examples of MNIST-M images misclassified by the baseline and correctly classified using SHADE, both trained with 250 samples.

adding SHADE, that encourages the model to discard information, allows important performance gains on this dataset and especially when only few training samples are given.

Finally, to confirm this behavior, we also applied the same procedure in a more conventional setting by training an Inception model on CIFAR-10. Figure 2b shows the results in that case. We can see that once again SHADE helps the model gain in performance and that this behavior is more noticeable when the number of samples is limited, allowing a gain of 6% when using 4000 samples.

## 4. CONCLUSION

In this paper, we introduced a new regularization method for DNNs training, SHADE, which focuses on minimizing the entropy of the representation conditionally to the labels. This regularization aims at increasing the intra-class invariance of the model while keeping class information. SHADE is tractable, adding only a small computational overhead when included into an efficient SGD training. We show that our SHADE regularization method significantly outperforms standard approaches such as weight decay or dropout with various DNN architectures on CIFAR-10. We also validate the scalability of SHADE by applying it on ImageNet. The invariance potential brought out by SHADE is further illustrated by its ability to ignore irrelevant visual information (texture, color) on MNIST-M. Finally, we also highlight the increasing benefit of our regularizer when the number of training examples becomes small.

## 5. REFERENCES

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NIPS*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., 2012.
- [2] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [3] Anders Krogh and John A. Hertz, "A simple weight decay can improve generalization," in *NIPS*, 1992.
- [4] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov, "Dropout: A simple way to prevent neural networks from overfitting," *JMLR*, 2014.
- [5] Sergey Ioffe and Christian Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *JMLR*, 2016.
- [6] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, "Understanding deep learning requires rethinking generalization," *ICLR*, 2017.
- [7] Stephane Mallat, "Group invariant scattering," in *Communications on Pure and Applied Mathematics*, 2012.
- [8] Joan Bruna and Stephane Mallat, "Invariant scattering convolution networks," in *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE*, 2013.
- [9] Gabriel Pereyra, George Tucker, Jan Chorowski, Lukasz Kaiser, and Geoffrey E. Hinton, "Regularizing neural networks by penalizing confident output distributions," in *ICLR Workshop*, 2017.
- [10] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," *Annual Allerton Conference on Communication, Control and Computing*, 1999.
- [11] T. Cover and J. Thomas, "Elements of information theory," *Wiley New York*, 1991.
- [12] A. A. Alemi, I. Fischer, J. V. Dillon, and K. Murphy, "Deep variational information bottleneck," in *ICLR*, 2017.
- [13] J. Rissanen, "Modeling by shortest data description," *Automatica*, 1978.
- [14] Tishby Naftali and Zaslavsky Noga, "Deep learning and the information bottleneck principle," in *Information Theory Workshop (ITW)*. IEEE, 2015.
- [15] Liam Paninski, "Estimation of entropy and mutual information," *Neural Computation*, 2003.
- [16] Diederik P. Kingma and Max Welling, "Auto-encoding variational bayes," in *ICLR*, 2014.
- [17] Xi Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, and Pieter Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *NIPS*, Jun 2016.
- [18] Gnter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter, "Self-normalizing neural networks," in *Advances in Neural Information Processing Systems 30*, Guyon I., Luxburg U. V., Bengio S., Wallach H., Fergus R., Vishwanathan S., and Garnett R., Eds., p. 972981. Curran Associates, Inc., 2017.
- [19] A. Krizhevsky, *Learning multiple layers of features from tiny images*, Ph.D. thesis, Computer Science Department University of Toronto, 2009.
- [20] Sergey Zagoruyko and Nikos Komodakis, "Wide residual networks," in *arXiv*, 2016.
- [21] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei, "Imagenet large scale visual recognition challeng," *ICJV*, 2015.
- [22] Durand Thibaut, Thome Nicolas, and Cord Matthieu, "WELDON: Weakly Supervised Learning of Deep Convolutional Neural Networks," in *CVPR*, 2016.
- [23] Ganin Yaroslav and Lempitsky Victor, "Unsupervised domain adaptation by backpropagation," in *ICML*, 2015, pp. 1180–1189.